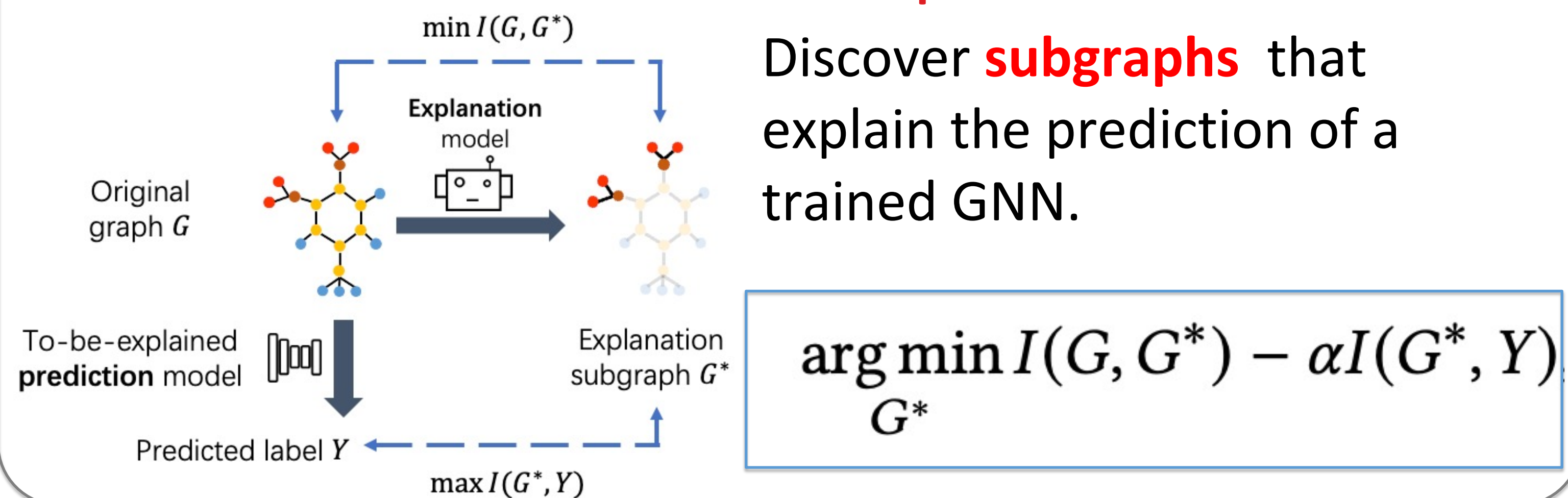# MixupExplainer: Generalizing Explanations for Graph Neural Networks with Data Augmentation

Jiaxing Zhang*, Dongsheng Luo*, Hua Wei

jz48@njit.edu, dluo@fiu.edu, hua.wei@asu.edu

**Arizona State University**

**FIU FLORIDA INTERNATIONAL UNIVERSITY**

**NJIT New Jersey Institute of Technology**

## Explaining GNNs

### Post-hoc instance-level explanation

$$\min I(G, G^*)$$

Discover **subgraphs** that explain the prediction of a trained GNN.

Original graph $G$

**Explanation model**

To-be-explained **prediction** model

Explanation subgraph $G^*$

Predicted label $Y$

$$\arg\min_{G^*} I(G, G^*) - \alpha I(G^*, Y)$$

$$\max I(G^*, Y)$$

## Graph Information Bottleneck (GIB) Objectives in a nut shell :
What was right and what was wrong?
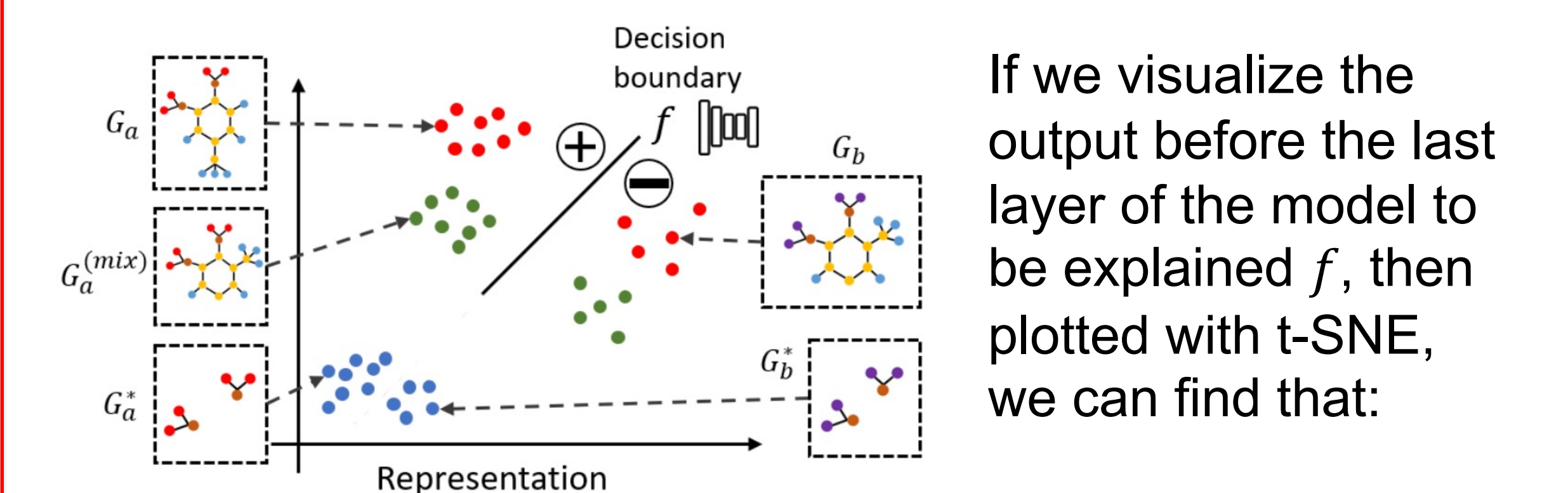
Mutual information $I(G^*, Y) = H(Y) - H(Y|G^*)$

*Intractability of $H(Y|G^*)$*
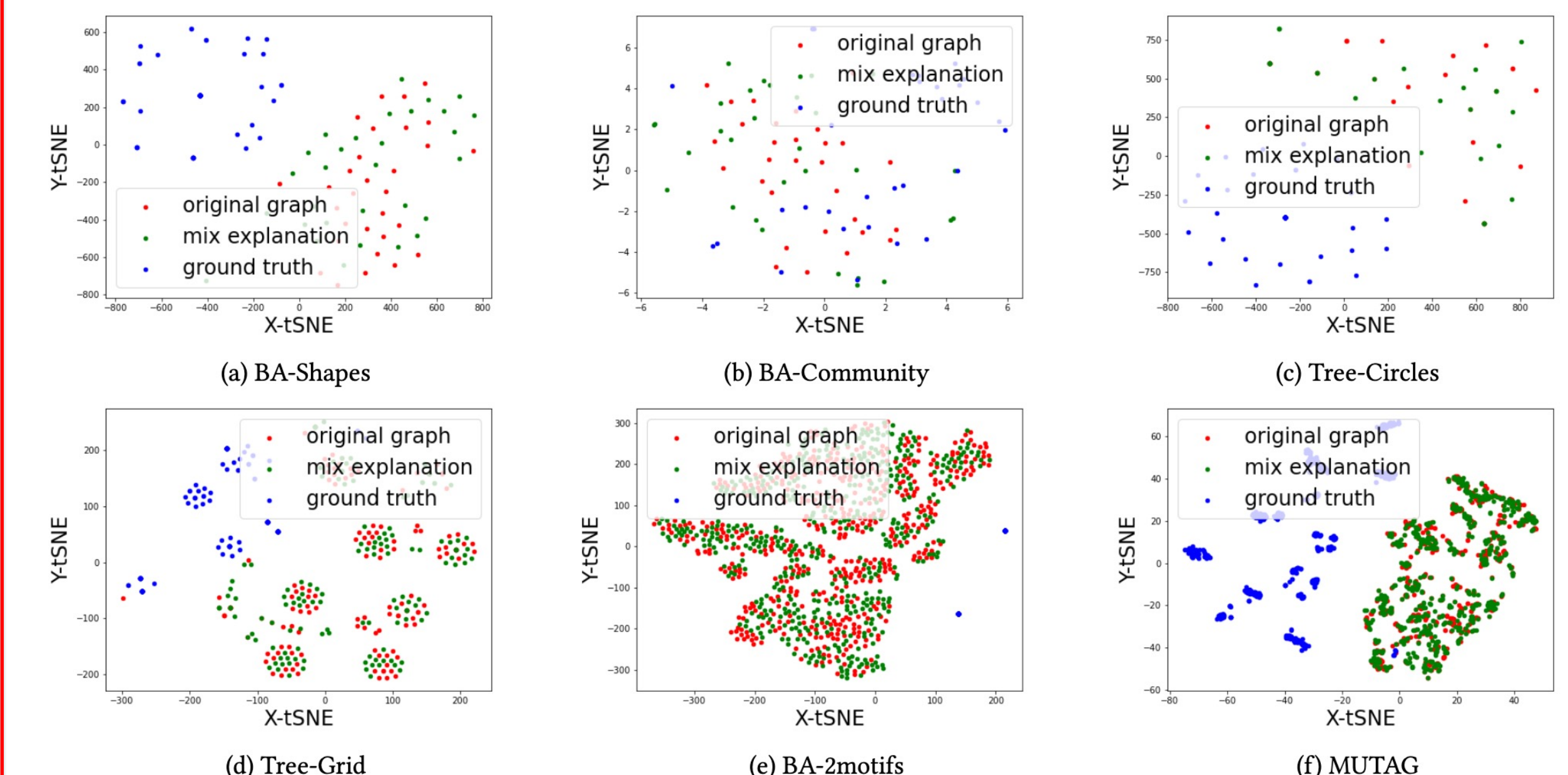
$$\arg\min_{G^*} I(G, G^*) + \alpha H(Y|G^*)$$

$$\arg\min_{G^*} I(G, G^*) + \alpha CE(Y, f(G^*))$$
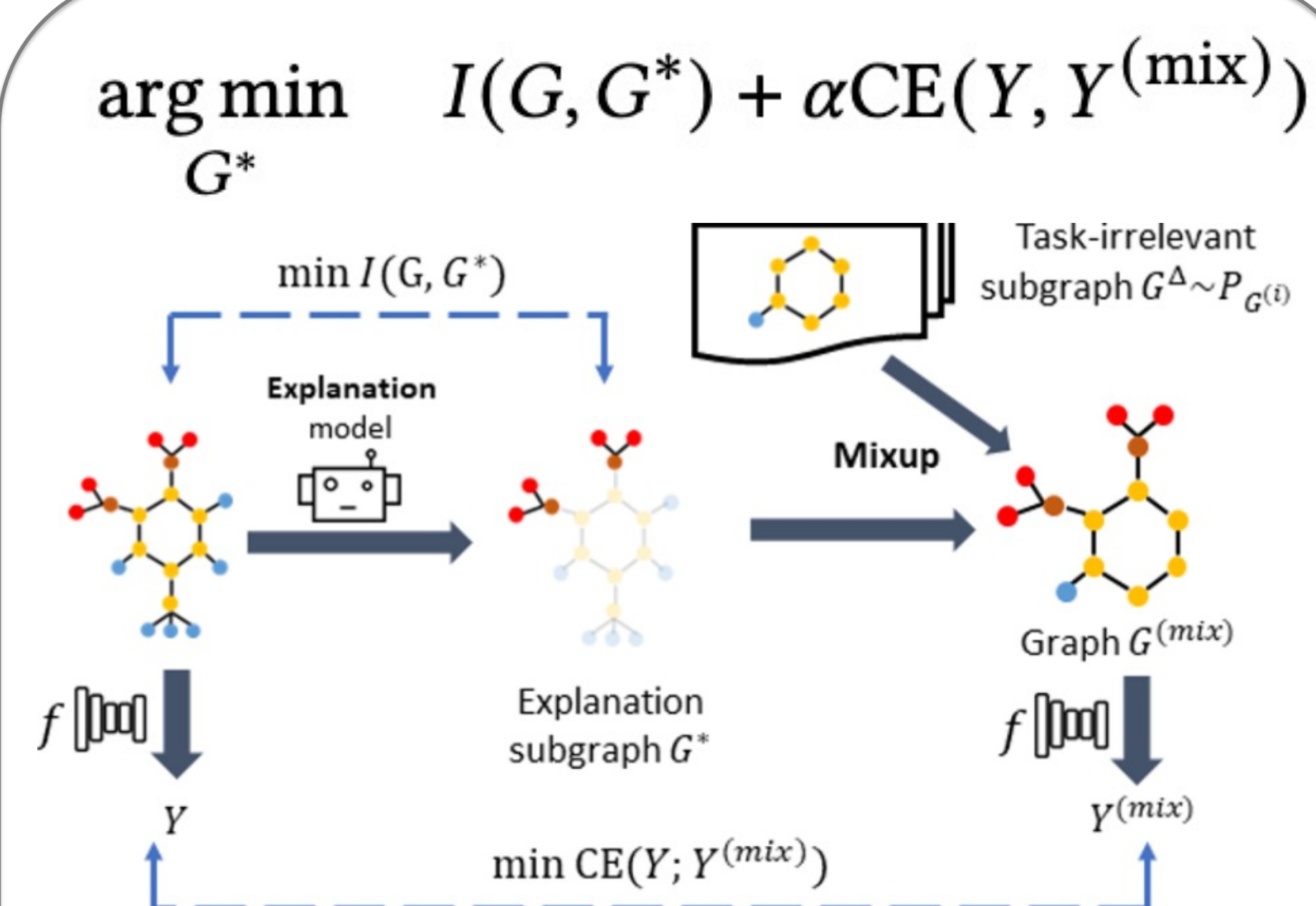[1,2,3]

### Diverging Distributions between $Y$ and $f(G^*)$

If we visualize the output before the last layer of the model to be explained $f$, then plotted with t-SNE, we can find that:

Decision boundary $f$

$G_a$

$G_a^{(mix)}$

$G_a$

$G_b$

$G_b^*$

Representation

**Ground truth** shifts away from **original graph**



(a) BA-Shapes    (b) BA-Community    (c) Tree-Circles

(d) Tree-Grid    (e) BA-2motifs    (f) MUTAG

**Mix explanation** aligns well with original graph

## Implementation

$$\arg\min_{G^*} \quad I(G, G^*) + \alpha CE(Y, Y^{(mix)})$$

$$\min I(G, G^*)$$

Task-irrelevant subgraph $G^\Delta \sim P_{G^{(t)}}$

**Explanation model**

**Mixup**

Explanation subgraph $G^*$

Graph $G^{(mix)}$

$f$ → $Y$

$f$ → $Y^{(mix)}$

$$\min CE(Y; Y^{(mix)})$$

**Intuition:** An explanation on a graph from GNN is the subgraph that can mix up with any random graphs and yet does not change GNN's prediction.

**Algorithm 1** Graph Mixup Algorithm

**Input:** Graph $G_a = (X_a, A_a)$, a set of graphs $\mathcal{G}$, the number of random connections $\eta$, explanation model $g$.

**Output:** Graph $G^{(mix)}$.

1: Randomly sample a graph $G_b = (A_b, X_b)$ from $\mathcal{G}$
2: Generate mask matrix $M_a = g(G_a)$
3: Generate mask matrix $M_b = g(G_b)$
4: Sample $\eta$ random connections between $G_a$ and $G_b$ as $A_c$
5: Mixup adjacency matrix $A^{(mix)}$ with Eq. (10)
6: Mixup edge mask $M^{(mix)}$ with Eq. (11)
7: Mixup node features $X^{(mix)} = [X_a; X_b]$
8: **return** $G^{(mix)} = (X^{(mix)}, M^{(mix)} \odot A^{(mix)})$

$$A^{(mix)} = \begin{bmatrix} A_a & A_c \\ A_c^T & A_b \end{bmatrix}$$

$$M_a^{(mix)} = \begin{bmatrix} \lambda M_a & M_c \\ M_c^T & A_b - \lambda M_b \end{bmatrix}$$

**Nice Property:** the proposed mixup approach could reduce the distance between the explanation and original graphs.

## Proposed generalized GIB objective

$$\arg\min_{G^*} I(G, G^*) + \alpha H(Y|G^*, G^\Delta) \quad \text{s.t.} \quad I(G^\Delta, Y|G^*) = 0.$$

**Nice property:** The generalized GIB objective is equivalent to vanilla GIB

$H(G)$ $H(Y)$

$G^*$ $S_1$ $S_2$

$G^\Delta$

$\boxed{}$ $I(G, G^*)$
$\boxed{}$ $S_1 = I(G^*, Y)$
$\boxed{}$ $S_2 = H(Y|G^*)$

$\boxed{}$ $S_3 = H(Y|G^*, G^\Delta)$
$\boxed{}$ $S_4 = I(G^\Delta, Y|G^*)$

When $S_4 = 0$, $S_3 = S_2$

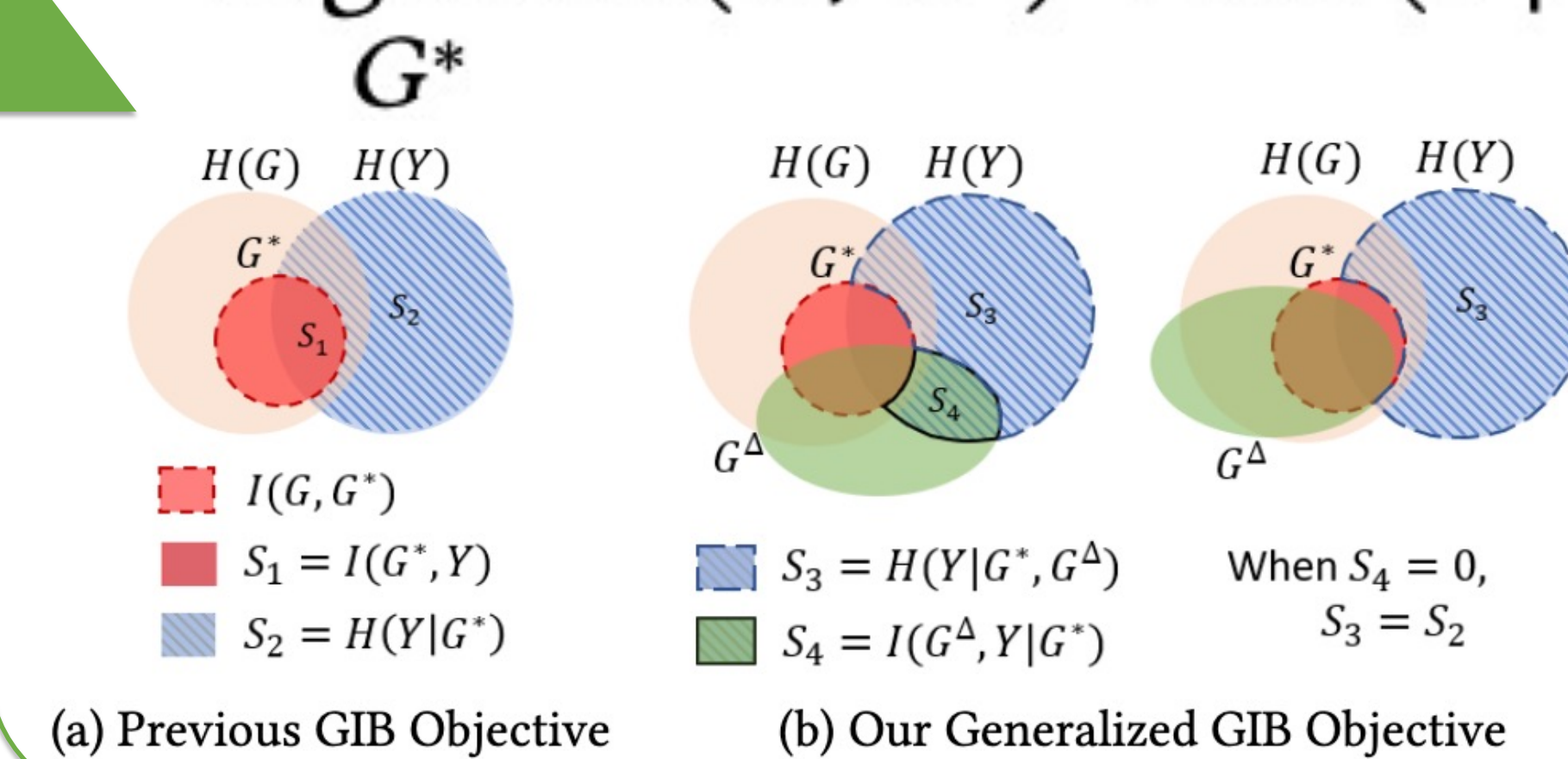(a) Previous GIB Objective    (b) Our Generalized GIB Objective

**Figure 2:** Illustration of GIB and our proposed objective. (a) Previous vanilla GIB objective aims to minimize $I(G^*, Y)$ and $H(Y|G^*)$, with a smaller overlap between $G^*$ and $G$. (b) Our generalized GIB objective has the same objective as vanilla GIB, with a larger lap between $G$ and $G^* + G^\Delta$, resulting in less distribution shifting issue.
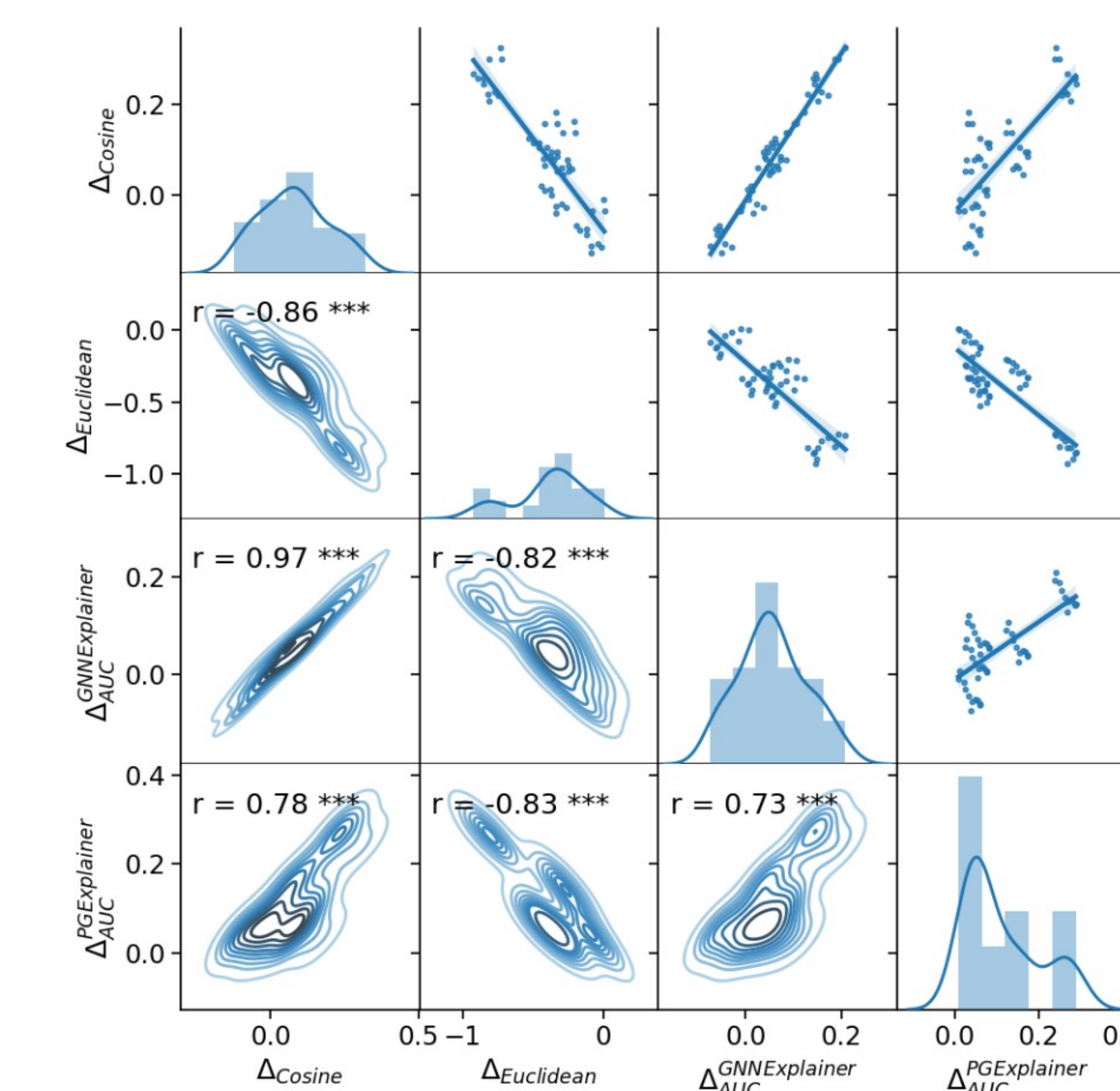
## Experiments

Table 1. Explanation faithfulness in terms of AUC-ROC on edges.

| | BA-Shapes | BA-Community | Tree-Circles | Tree-Grid | BA-2motifs | MUTAG |
|---|---|---|---|---|---|---|
| GRAD | 0.882 | 0.750 | 0.905 | 0.612 | 0.717 | 0.783 |
| ATT | 0.815 | 0.739 | 0.824 | 0.667 | 0.667 | 0.765 |
| SubgraphX | 0.548 | 0.473 | 0.617 | 0.516 | 0.610 | 0.529 |
| MetaGNN | 0.851 | 0.688 | 0.523 | 0.628 | 0.500 | 0.680 |
| RG-Explainer | 0.985 | 0.919 | 0.787 | 0.927 | 0.657 | 0.873 |
| GNNExplainer | $0.884_{\pm0.002}$ | $0.682_{\pm0.004}$ | $0.683_{\pm0.009}$ | $0.379_{\pm0.001}$ | $0.660_{\pm0.006}$ | $0.539_{\pm0.002}$ |
| + MixUp | $0.890_{\pm0.004}$ | $0.788_{\pm0.006}$ | $0.690_{\pm0.014}$ | $0.501_{\pm0.003}$ | $0.869_{\pm0.004}$ | $0.612_{\pm0.043}$ |
| (improvement) | 0.60% | 15.5% | 1.02% | 32.2% | 31.7% | 13.5% |
| PGExplainer | $0.999_{\pm0.001}$ | $0.829_{\pm0.040}$ | $0.762_{\pm0.014}$ | $0.679_{\pm0.008}$ | $0.679_{\pm0.043}$ | $0.843_{\pm0.084}$ |
| + MixUp | $0.999_{\pm0.001}$ | $0.955_{\pm0.017}$ | $0.774_{\pm0.004}$ | $0.712_{\pm0.000}$ | $0.920_{\pm0.031}$ | $0.871_{\pm0.079}$ |
| (improvement) | 0.00% | 15.2% | 1.57% | 4.86% | 35.5% | 3.32% |



All these four improvements strongly correlated to each other with statistical significance: The improvements achieved by MixupExplainer in explanations accuracy own to the successful alleviation of the distribution shifting issue.

## References

[1]. Ying, et al, GNNexplainer: Generating explanations for graph neural networks. NeurIPS 2020.
[2]. Luo, et al., Parameterized explainer for graph neural network. NeurIPS 2020.
[3]. Miao et al., Interpretable and generalizable graph learning via stochastic attention mechanism. ICML 2022.

## Conclusion

- **Be careful** if you are using GNNExplainer or PGExplainer! You might encounter distribution shifting issue between $f(G)$ and $f(G^*)$.
- An explanation of a GNN's prediction on an original graph is the subgraph that can mix up with any random graphs and does not change GNN's prediction.